



A software tool for the spatiotemporal analysis and reporting of groundwater monitoring data



Wayne R. Jones^{a,*}, Michael J. Spence^a, Adrian W. Bowman^b, Ludger Evers^b,
Daniel A. Molinari^b

^a Shell Global Solutions (UK), Brabazon House, Concord Business Park, Threapwood Road, Manchester, M22 0RR, United Kingdom

^b School of Mathematics and Statistics, University Gardens, University of Glasgow, Glasgow G12 8QQ, United Kingdom

ARTICLE INFO

Article history:

Received 22 May 2013

Received in revised form

15 January 2014

Accepted 16 January 2014

Available online 25 February 2014

Keywords:

Environmental monitoring

Groundwater

Open source software

R

Spatiotemporal

Geostatistics

Spatial modelling

GWSDAT

ABSTRACT

The GroundWater Spatiotemporal Data Analysis Tool (GWSDAT) is a user friendly, open source, decision support tool for the analysis and reporting of groundwater monitoring data. Uniquely, GWSDAT applies a spatiotemporal model smoother for a more coherent and smooth interpretation of the interaction in spatial and time-series components of groundwater solute concentrations. Data entry is via a standardised Microsoft Excel input template whilst the underlying statistical modelling and graphical output are generated using the open source statistical program R. This paper describes in detail the various plotting options available and how the graphical user interface can be used for rapid, rigorous and interactive trend analysis with facilitated report generation. GWSDAT has been used extensively in the assessment of soil and groundwater conditions at Shell's downstream assets and the discussion section describes the benefits of its applied use. Finally, some consideration is given to possible future developments.

© 2014 The Authors. Published by Elsevier Ltd. Open access under [CC BY license](http://creativecommons.org/licenses/by/4.0/).

Software availability

Software name: GWSDAT (GroundWater Spatiotemporal Data Analysis Tool)

Developer: Wayne R. Jones

Contact address: Shell Global Solutions (UK)

(wayne.w.jones@shell.com)

Year first official release: 2013

Hardware requirements: Standard PC

System requirements: Microsoft Windows (XP or later)

Software requirements: Microsoft Office (Excel, Word and PowerPoint) and R (www.r-project.org)

Program Size: 13 MB

Availability: www.claire.co.uk/GWSDAT

License: Free under a GNU General Public License (www.gnu.org) agreement.

Documentation and support for users: User manual, example data sets, FAQ document, presentations and posters.

1. Introduction

1.1. Background

Groundwater is water located beneath the Earth's surface in soil pore spaces and in the fractures of rock formations. Environmental monitoring of groundwater is routinely conducted in areas where the risk of contamination is high and for protecting human health and the environment following an accidental release of hazardous constituents. Groundwater monitoring strategies are designed to establish the current status and assess trends in environmental parameters, and to enable an estimate of the risks to human health and the environment. It involves installing a network of monitoring wells to enable access to the water table across the site (Barcelona et al., 1985). Samples of groundwater are periodically collected from these wells and sent to an accredited laboratory for chemical analysis. The resulting spatiotemporal data set has to be reviewed,

* Corresponding author.

E-mail address: wayne.w.jones@shell.com (W.R. Jones).

GWSDAT (GroundWater Spatio-Temporal Data Analysis Tool)
 Author: Wayne.W.Jones@Shell.com Version: 2.0

Historical Monitoring Data					
WellName	Constituent	SampleDate	Result	Units	Flags
MW-01	BENZENE	31/10/2002	78	ug/l	
MW-01	GW	31/10/2002	92.23	Level	
MW-01	TOLUENE	31/10/2002	470	ug/l	
MW-01	XYLENE	31/10/2002	430	ug/l	
MW-02	BENZENE	31/10/2002	40000	ug/l	
MW-02	GW	31/10/2002	92.3	Level	
MW-02	TOLUENE	31/10/2002	1200	ug/l	
MW-02	XYLENE	31/10/2002	1100	ug/l	
MW-03	BENZENE	31/10/2002	ND<10	ug/l	
MW-03	GW	31/10/2002	94.43	Level	
MW-03	TOLUENE	31/10/2002	1.1	ug/l	
MW-03	XYLENE	31/10/2002	ND<0.50	ug/l	
MW-04	BENZENE	31/10/2002	250	ug/l	

Well Coordinates			
WellName	XCoord	YCoord	Aquifer
MW-01	97.43	57.81	
MW-02	85.57	50.64	
MW-03	22.95	74.64	
MW-04	83.64	81.26	
MW-05	42.26	114.64	
MW-06	62.40	44.57	
MW-07	126.12	72.43	
MW-08	126.95	104.15	
MW-09	141.84	42.09	
MW-10	111.50	23.05	
MW-11	88.05	7.88	

GIS ShapeFiles	
FileNames (*.shp)	

Fig. 1. GWSDAT example data input template. The *Historical Monitoring Data* table captures the concentration data, groundwater levels and, if present, NAPL thickness. The *Well Coordinates Table* stores the location of the monitoring well. The GWSDAT add-in menu is displayed at the top left.

analysed statistically, interpreted, and the results presented to environmental regulators in a clear and understandable manner.

The most basic method of level and trend evaluation involves investigating the time-series of groundwater constituent concentrations independently on a well by well basis. The more sophisticated spatial methods, typically, involve fitting a concentration trend surface (i.e. Kriging) to evaluate spatial pattern and trend (Cameron and Hunter, 2002; Gaus et al., 2003). However, although spatiotemporal data lies at the heart of current research in statistical methods (see Cressie and Wikle (2011)), the most common practice is to independently apply spatial modelling techniques to separate monitoring events (e.g. Ricker (2008)) or apply a single spatial model to a data set which has been consolidated over a time period (e.g. Aziz et al. (2003)). The joint modelling of both spatial and time elements in a single spatiotemporal modelling framework leads to a more coherent interpretation of site groundwater characteristics (Evers et al., in press).

Whilst there is a range of freely available groundwater data analysis applications, the most sophisticated tend to be designed for large scale long term groundwater monitoring networks (Aziz et al., 2003; Cameron, 2004). These have a relatively large initial data warehousing setup burden, which may be viewed as a barrier to the more widespread use of advanced groundwater monitoring techniques to smaller more short term monitoring programmes. Similarly, whilst GIS applications (e.g. ArcGIS) have excellent visualisation tools for geographical interpretation they also have a high initial data setup cost, operator competence requirements, and perhaps surprisingly, only a limited number of geostatistical modelling techniques available.

2. Software design and aims

2.1. Development aims

To a large extent, GWSDAT has been developed to address the barriers discussed in Section 1.1. However, its most important aim is to provide a simple to use, but statistically powerful decision support tool to environmental engineers and practitioners who routinely report on the status of numerous groundwater monitoring sites. Such an application needs to be easy to setup yet

powerful in its ability to objectively analyse and rapidly report on a groundwater monitoring site's characteristics.

In common with many other environmental applications, it was recognised that there would be a benefit in providing the software in an open and transparent manner because policy makers and environmental regulators generally prefer code and techniques which are fully transparent and supported by sound science (Carslaw and Ropkins, 2012).

2.2. Software architecture

GWSDAT has been designed to integrate with Microsoft Excel, a software routinely used by environmental engineers for storing and analysing environmental (e.g. soil and groundwater) data. The user entry point to GWSDAT is a custom built Excel Add-in menu (see top left of Fig. 1).

The statistical engine used to perform geostatistical modelling and display graphical output is the open source statistical programming language R (R Development Core Team, 2012). The R project is used across a wide range of disciplines and has been adopted with eagerness by the environmental sciences community (Carslaw and Ropkins, 2012). Members of the R community contribute statistical routines and functionality to this collaborative project by means of an open standardised package structure, which can be downloaded and installed from <http://cran.r-project.org/web/packages/>. GWSDAT makes use of several of these packages, which are all individually referenced in this article. A Graphical User Interface (GUI) is provided via the R packages *rpanel* (Bowman et al., 2007) and *tkrplot* (Tierney, 2011) which obviates the need for training GWSDAT users in the R programming language.

3. Data input

3.1. Background

Before describing the application of GWSDAT in more detail it is necessary to give a brief explanation of the nature of groundwater monitoring data. In general, routine sampling of a monitoring well involves measuring the groundwater elevation and taking a sample of the groundwater which is subsequently sent for laboratory

analysis to ascertain the dissolved concentration of a prescribed set of solutes (e.g. Toluene, Benzene). If the concentration is deemed lower than that which could be detected using the method employed by the laboratory then it is classified as a 'non-detect'. In such circumstances, the laboratory quotes the detection threshold concentration value below which the solute could not be detected.

An additional important consideration for petroleum hydrocarbon applications is the presence of a layer of Non-Aqueous Phase Liquid (NAPL), such as gasoline or diesel, on the surface of the water table. This circumstance often arises when the amount of contamination is sufficient to exceed the natural solute level of groundwater. Samples containing NAPL are not often sent for a full chemical analysis (unless performing NAPL forensics) because the levels of solute concentrations are too high for the traditional laboratory methods, which are geared towards lower concentrations. Hence, NAPL data poses the challenge of how to handle unspecified high solute concentration values and identify trends in NAPL layer thickness.

3.2. Input data format

Groundwater monitoring data is entered into GWSDAT by means of a simple standardised Microsoft Excel input sheet (Fig. 1). There is no requirement to gather any data that would not have already been recorded in a standard groundwater monitoring data set. The following summarises the GWSDAT data input format but the reader is referred to the user manual for a full and detailed explanation of GWSDAT data input specification.

Each row of the *Historical Monitoring Data* table (left hand table in Fig. 1) corresponds to a unique combination of well id, sampling date, aquifer zone, solute name and concentration. Non-detect solute data is entered using the notation 'ND < X', where X represents the laboratory reported detection threshold concentration. If present, NAPL thickness data is also entered in this table using the constituent name 'NAPL' with an appropriate unit, e.g. metres, mm. Optionally, groundwater level data is entered here (using the constituent name 'GW') as an elevation above a common datum, e.g. metres or feet above sea level or some other common reference height.

The *Well Coordinates* table (middle table in Fig. 1) stores the coordinates of the groundwater monitoring wells. Any arbitrary coordinate system with an aspect ratio of 1 can be used, i.e. a unit in the x-coordinate is the same distance as a unit in the y-coordinate.

The third optional *GIS Shapefiles* table can be populated with file locations of GIS shapefiles (Esri, 1998) for use as basemaps or site plans. Two GWSDAT input data sets of varying complexity (basic and comprehensive) are included with the software to serve as both an example of the GWSDAT data input format and provide a quick way of getting started.

3.3. Data processing

On initiation of a GWSDAT analysis, the user is asked to select from a variety of data processing options including the handling of non-detects and, if present, NAPL. In accordance with the common convention, the default option is to substitute the non-detect solute concentration data with half its detection limit. Note, however, that this can mask trends in the data and lead to erroneous estimates of summary statistics in cases with a high proportion of non-detect results (Helsel, 2004). This issue is discussed in more detail in Section 6. If NAPL is present the user is prompted to substitute NAPL data points with site data set maximum observed solute concentrations. This option is to provide a more realistic picture of the area of impacted groundwater (high concentrations) in the event that NAPL in wells prevents direct measurement of

solute concentrations as discussed in Section 3.1. The data processing step is concluded with a series of data validation procedures to check for common data input errors.

4. Graphical user interface

4.1. Introduction

In the interests of user-friendliness and productivity the results of a GWSDAT analysis are interrogated and interpreted through the GWSDAT user interface (see Fig. 2). It includes a wide range of different plots for the visual inspection of groundwater monitoring data. The objective assessment of trend is achieved by the application of statistical smoothing models described in Appendix A. The following sections describe the individual components of the GWSDAT user interface in more detail.

4.2. Well trend plot

The well trend plot (see Fig. 3) enables the user to investigate time series trends of solute concentrations and groundwater level in individual wells. Sampled concentration values are displayed using orange circles for non-detect data and black solid circles for detectable data. The user can choose to overlay a linear (or log-linear) regression model fit and use the non-parametric Mann–Kendall approach to trend detection via the R package *Kendall* (McLeod, 2011). Although this approach is widely used in environmental sciences (e.g. Hirsch et al. (1982); Helsel and Hirsch (2002)) its major weakness is that it can only detect monotonic trend and in response GWSDAT adopts an additional methodology. The solid blue line in Fig. 3 displays the estimate (together with a 95% confidence interval) of the mean trend level according to a local linear regression model fit described in Appendix A.1. This non-parametric model smoothing technique is not constrained to be monotonic and can change direction as is clearly illustrated in the figure. The trend between two points in time is, informally speaking, deemed statistically significant if the associated confidence intervals do not overlap (Fig. 4).

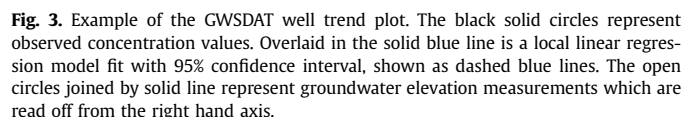
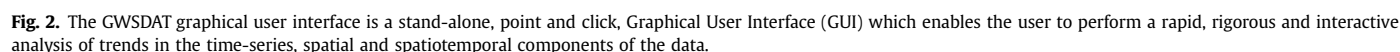
For evaluating the impact of changing (perhaps seasonal) water table conditions groundwater elevation data can, optionally, be overlaid in this plot. The time series of observed groundwater level is represented by open circles joined by a black solid line see and the values read off from the right hand axis (see Fig. 3). If present, NAPL thickness data can also be displayed in a similar manner.

4.3. Trend and threshold indicator matrix

The trend and threshold indicator matrix is a summary of the level and time-series trend in solute concentrations at a particular time-slice of the monitoring period. The rows correspond to each monitoring well and the columns correspond to the different solutes. The date of the time-slice is displayed at the top of the plot and also indicated by a vertical grey line in the well trend plot (see Fig. 3). The user can select between the options of displaying 'Trend', 'Threshold – Absolute' or 'Threshold Statistical'.

When 'Trend' is selected the cells are coloured to indicate the strength and direction of the current trend as assessed by the instantaneous gradient of the well trend smoother (see Section 4.2) at the current time-slice. White cells indicate a generally flat trend whilst reds and greens indicate strong upward and downward trends, respectively. In the event that the trend cannot be calculated (e.g. no data) then the corresponding cell is coloured grey. Blue cells represent non-detect data.

When 'Threshold Absolute' is selected the cells are coloured according to whether the observed current solute concentrations



are below a user specified threshold value, such as a risk-based remedial objective. The cells are coloured red if the current solute concentration is above the threshold value and green otherwise. 'Threshold Statistical' is similar but only colours the cell green if the upper 95% confidence interval of the well trend smoother (see Section 4.2) is below the threshold value.

4.4. Spatial plot

The GWSDAT spatial plot (see Fig. 5) is for the analysis of spatial trends in solute concentrations, groundwater flow and, if present, NAPL thickness. It displays the locations of the named monitoring wells together with sample solute concentration values collected within the date interval displayed at the top of the graphic. If desired, the major site features (e.g. roads, fuel tanks), supplied in a GIS shapefile format, can be overlaid on the spatial plot as light blue lines. As the user increments forwards and backwards through the monitoring history, using the ‘+’ and ‘-’ *Time Steps* buttons, the spatial plot is updated.

The estimated groundwater flow direction and magnitude is depicted with blue arrows calculated using the method described in [Appendix A.2](#). Additionally, it is possible to overlay a contour plot of groundwater elevation. This is achieved by drawing isopleths through a fitted local polynomial regression model fit implemented

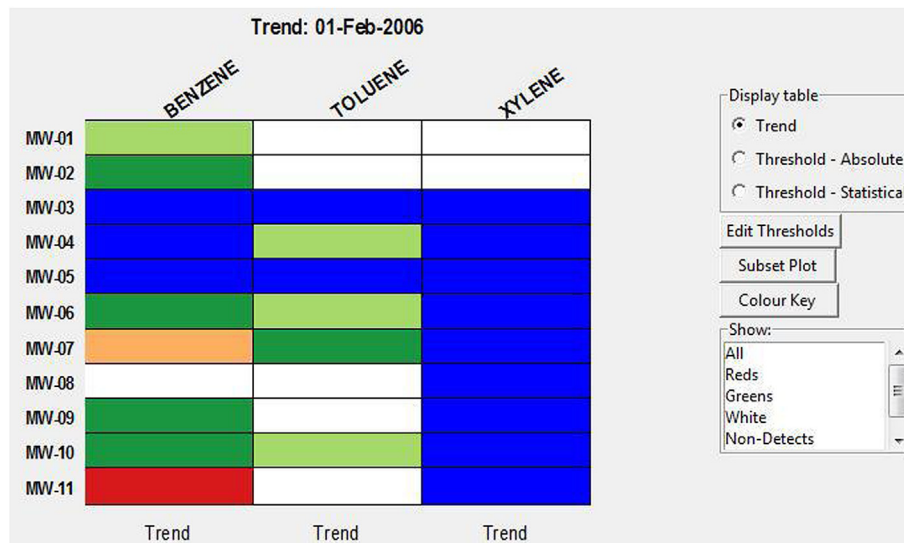


Fig. 4. Example of the GWSDAT trend and threshold indicator matrix. The rows represent monitoring wells, and columns represent the different solutes. Each cell is colour coded to represent increasing (reds), stable (white) or decreasing (greens) trends in solute concentrations. Blue cells represent non-detect data and if there is insufficient data the cell is coloured grey.

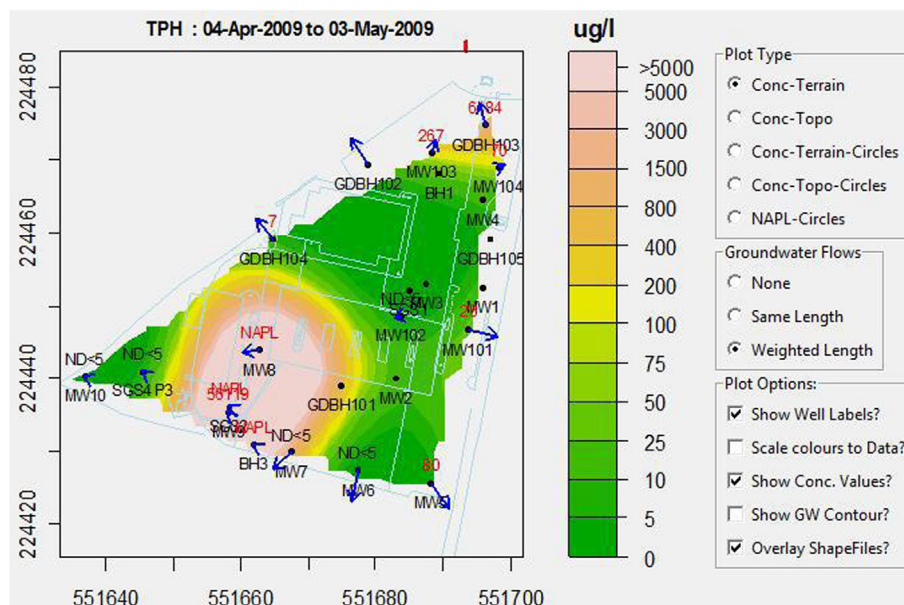


Fig. 5. Example of the GWSDAT spatial plot. The location of the named monitoring wells is depicted with black solid circles. Detect data or NAPL is displayed in a red font and non-detect in a black font above the wells. Blue arrows indicate vectors of estimated groundwater flow velocity. Spatiotemporal solute concentration smoother predictions are superposed using the colour key on the right. GIS shapefile data is overlaid using light blue lines.

using the R function *loess* – a 2D variant of the local linear regression method explained in [Appendix A.1](#).

The spatial distribution of solute concentration is estimated by taking a time-slice through the spatiotemporal concentration smoother (discussed further in [Section 4.5](#)). The model predictions are superposed on the spatial plot with a user-specified colour key located to the right of the plot. Alternatively, if no model based predictions are required, the concentration smoother can be replaced by size scaled colour coded circles representing the magnitude of sampled solute concentration values. If NAPL is present, the additional 'NAPL-Circles' option is available which displays NAPL thickness measurements at the monitoring well locations using a similar circle based representation, i.e. a bubble-plot.

The spatial plot uses the R packages, *sp* ([Pebesma and Bivand, 2005](#)), *splancs* ([Rowlingson et al., 2012](#)) and *maptools* ([Lewin-Koh et al., 2012](#)).

4.5. Spatiotemporal trend analysis

One of GWSDAT's unique features is that the spatial and temporal components of the solute concentration data are modelled jointly in a single modelling framework described in [Appendix A.3](#). The simultaneous statistical smoothing of both spatial and temporal components provides a clearer and more insightful interpretation of the groundwater monitoring site solute characteristics than would otherwise be gleaned from analysing these two components

separately. However, it is not an inconsiderable challenge to effectively communicate the 3-dimensional nature of spatiotemporal trend through a 2-dimensional medium of a computer monitor. Furthermore, there is an additional constraint that the output from a GWSDAT analysis is commonly used in paper-based non-interactive reports submitted to environmental regulators. For this reason, GWSDAT communicates spatiotemporal trend through automatic plotting of the full temporal sequence of spatial plots (see Section 4.4). This animation based approach provides a ‘movie’ clearly illustrating how both the spatial and temporal distribution of historical groundwater solute concentrations have changed over the monitoring period.

The ‘animations’ menu located at the top-left of the GWSDAT user interface (Fig. 2) provides three different methods for generating animations. The first method plots and records the full sequence of spatial plots in an R graphics window. The user can toggle forwards and backwards through the sequence of spatial plots using the ‘Page Up’ and ‘Page Down’ keyboard buttons. The second method is identical but additionally generates a Microsoft PowerPoint slide-pack of the full sequence of spatial plots. The third method uses the R package *animation* (Xie, 2012), to generate a html animation page (with controls) of spatial plots in the user’s internet browser. The html animation can be viewed independently of GWSDAT, and hence provides an excellent dynamic media for communicating results to individuals who do not have direct access to GWSDAT.

4.6. Report generation

By left-clicking on any of the GWSDAT user interface plots, an identical but expanded plot is generated in a separate R graphics window. Plots can be saved to a variety of different formats including ‘jpeg’, ‘postscript’, ‘pdf’, ‘metafile’. Alternatively, with a

single click of a mouse, plots and sequences of plots (e.g. spatio-temporal animation described in Section 4.5) can be diverted directly in to Microsoft Word or PowerPoint. This functionality, implemented using the R package *RDCOMClient* (Lang, 2012), enables the user to interactively compile a site groundwater monitoring report in an expeditious manner.

Additional report generation functionality include the ‘Well Reporting’ procedure, implemented using the R package *lattice* (Sarkar, 2008), which generates a matrix of graphs displaying time series solute concentration values on a well by well basis (see Fig. 6). This plot can be used to very concisely summarise the time series trends in the complete set of solutes and monitoring wells. A similar report procedure ‘GW Well Reporting’ also allows for the overlay of the time series in groundwater elevation at each well. Finally, the ‘Latest Snapshot’ procedure generates a sequence of plots (to PowerPoint if required) which reports on the most recent trends. This includes the latest spatial plot for each solute together with the most recent three variants of the ‘Trend and Threshold Indicator Matrix’ plot described in Section 4.3.

5. Discussion

Environmental risk-based management decisions are often based on limited understanding of groundwater data, and relatively limited statistical analysis of that data. GWSDAT has been designed and developed as a user-friendly, interactive, trend analysis tool for distilling the information from such groundwater monitoring data sets. The application has been used operationally in the monitoring and assessment of Shell’s global downstream assets (e.g. refineries, terminals, fuel stations) for a period of over 4 years. Graphical output generated from GWSDAT is routinely included in reports

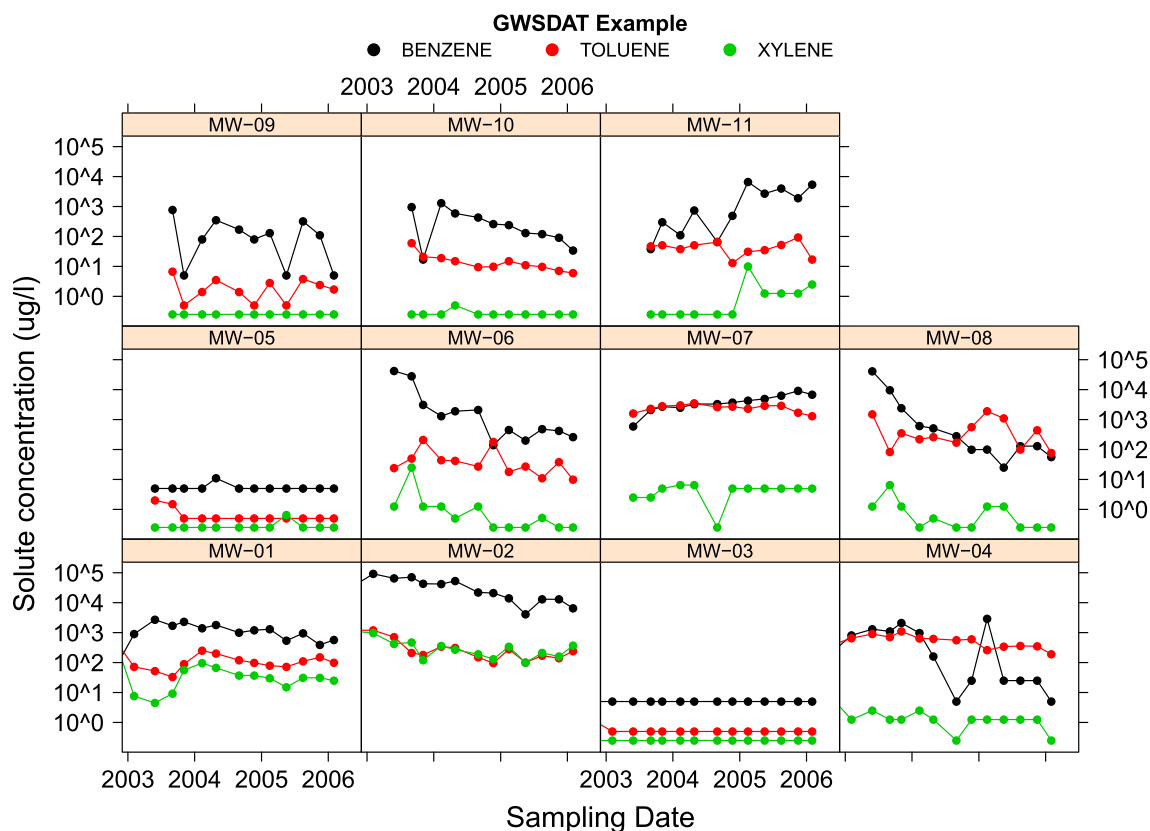


Fig. 6. Example of the GWSDAT Well report plot. The colour key at the top identifies each solute and the name of each well is displayed in a banner at the top of each of the individual time series graphs. This clearly illustrates the correlation in time series trends across the different solutes.

submitted to environmental regulators. Environmental engineers using GWSDAT have reported numerous benefits:

- Rapid interpretation of complex data sets for both small and large groundwater monitoring networks.
- Earlier identification of new spills or off-site migration.
- Reduced reliance on engineered remediation through increased use of monitored natural attenuation remedies, where groundwater data analysis supports its effectiveness.
- Earlier closeout of sites in needless long-term monitoring and/or remediation.
- Simplified preparation of groundwater monitoring reports.

6. Future developments

The major area for future development is the addition of new capabilities to GWSDAT. The assessment of solute plume stability is currently carried out by visually inspecting the evolution of the spatiotemporal solute concentration smoother. Feedback from users has highlighted the need for additional quantitative tools to supplement this graphical method. Development is currently underway to incorporate plume mass balance tools, such as those proposed in [Ricker \(2008\)](#), to automatically estimate plume characteristics such as area, total mass and centre of mass. The inspection of these quantities over the monitoring period will more objectively illustrate whether the plume is moving and if the plume is growing, shrinking or stable.

Future versions of GWSDAT will use spatiotemporal model standard errors to give the user a better understanding of model uncertainty and goodness of fit. The spatial distribution of model standard errors is of particular interest because it provides an assessment of the design of the well monitoring network. Areas of low monitoring density will have larger model standard errors. This not only informs the user that the predictions in this area need to be interpreted with care but also identifies potential locations where the construction of new monitoring wells would improve conceptual understanding of a site, and project decision-making. Model standard errors could also be used in the calculation of the solute plume characteristics mentioned above to provide a confidence interval on these quantities.

Whilst simple to implement, the substitution of non-detect concentration values is not without its disadvantages as discussed by [Helsel \(2004\)](#). These are partly mitigated in GWSDAT by offering the 'worse case' scenario of substitution with the full detection limit as opposed to the usual value of half the detection limit. However, the occurrence of different detection limits for the same solute (perhaps because different laboratories were used during the course of a long-term monitoring programme) is still troublesome as substitution with any constant fraction leads to an apparent trend in concentrations. The authors are currently researching more sophisticated censored regression techniques to handle non-detect data in the spatiotemporal modelling framework.

Acknowledgements

This work was funded by Shell Global Solutions (UK) Ltd. The authors acknowledge contributions from numerous colleagues to the development of GWSDAT: Dr Matthew Lahvis, Dr George Devaull, Dan Walsh, Curtis Stanley, and Professor Jonathan Smith of Shell Projects & Technology HSE Technology; Philip Jonathan of Projects & Technology – Analytical Services: Statistics & Chemometrics; Ewan Crawford, of Glasgow University, Scotland, UK. The views expressed are those of the authors and may not reflect the policy or position of Royal Dutch Shell plc.

Appendix A. Description of statistical modelling techniques

Appendix A.1. Well trend plot smoother

The well trend plot smoother is fitted using a non-parametric method called local linear regression. This involves solving locally the least squares problem:

$$\min_{\alpha, \beta} \sum_i^n \{y_i - \alpha - \beta(x_i - x)\}^2 w(x_i - x; h) \quad (\text{A.1})$$

where $w(x_i - x; h)$ is a weight function with parameter h . The weight function gives the most weight to the data points nearest the point of estimation and the least weight to the data points that are furthest away. For the weight function GWSDAT uses a normally-distributed probability density function with standard deviation h . Local linear regression is deployed in GWSDAT using the R package *sm* ([Bowman and Azzalini, 2010, 1997](#)) and h is selected using the method published in [Hurvich et al. \(1998\)](#).

Appendix A.2. Groundwater flow estimation

Vectors of groundwater flow strength and direction are estimated using the well coordinates and recorded groundwater elevations. The model is based on the simple premise that local groundwater flow will follow the local direction of steepest descent (hydraulic gradient).

For a given well, a linear plane is fitted to the local groundwater level data:

$$L_i = a + bx_i + cy_i + \epsilon_i \quad (\text{A.2})$$

where L_i represents the groundwater level at location (x_i, y_i) . Local data is defined as the neighbouring wells as given by a Delaunay triangulation ([Ahuja and Schacter, 1983; Turner, 2012](#)) of the monitoring well locations. The gradient of this linear surface in both x and y directions is given by the coefficients b and c . Estimated direction of flow is given by:

$$\theta = \tan^{-1}\left(\frac{c}{b}\right) \quad (\text{A.3})$$

and the relative hydraulic gradient (a measure of relative flow velocity) is given by

$$R = \sqrt{b^2 + c^2} \quad (\text{A.4})$$

For any given model output interval this algorithm is applied to each and every well where a groundwater elevation has been recorded.

Appendix A.3. Spatiotemporal solute concentration smoother

The spatiotemporal solute concentration smoother is fitted using a non-parametric regression technique known as Penalised Splines ([Eilers and Marx, 1992, 1996](#)). A full and detailed explanation of applying this statistical method to groundwater monitoring data is the subject of another paper ([Evers et al., 2014](#)). However, the following outlines some of the most important aspects for the purposes of GWSDAT.

Let y_i be the natural log solute concentration at $x_i = (x_{i1}, x_{i2}, x_{i3})$ where x_{i1} and x_{i2} stand for the spatial coordinates of the well and x_{i3} represents the corresponding time point for the i -th observation with $i = 1, \dots, n$. We start by modelling the solute concentration as

$$y_i = \sum_{j=1}^m b_j(\mathbf{x}_i) \alpha_j + \epsilon_i \quad (\text{A.5})$$

where the $b_j, j = 1, \dots, m$ are m B-Spline basis functions, generally second or third order polynomials (Eilers and Marx, 1996). The measurement errors ϵ_i 's are assumed to be independent and identically normally distributed with zero mean and variance σ^2 . Rewriting Equation (A.5) in the more compact matrix notation leads to

$$\mathbf{y} = \mathbf{B}(\mathbf{x})\boldsymbol{\alpha} + \boldsymbol{\epsilon} \quad (\text{A.6})$$

The traditional ordinary least squares approach is to minimize the objective function $S(\boldsymbol{\alpha}) = \|\boldsymbol{\epsilon}\|^2 = \|\mathbf{y} - \mathbf{B}(\mathbf{x})\boldsymbol{\alpha}\|^2$. The well known major disadvantage of this approach is its propensity to overfit data leading to under smoothness in model predictions. To overcome this hurdle, the objective function is modified with the addition of a term that penalises the finite differences of the coefficients of adjacent B-splines. The objective function now takes the form $S(\boldsymbol{\alpha}) = \|\mathbf{y} - \mathbf{B}(\mathbf{x})\boldsymbol{\alpha}\|^2 + \lambda\|D_d\boldsymbol{\alpha}\|^2$ where D_d is a matrix such that $D_d = \Delta^d$, the d -th differences of $\boldsymbol{\alpha}$, and λ is a nonnegative tuning parameter.

By minimising the new objective function for a given value of λ , we obtain the estimator of the parameters $\hat{\boldsymbol{\alpha}} = (\mathbf{B}^T\mathbf{B} + \lambda\mathbf{D}_d^T\mathbf{D}_d)^{-1}\mathbf{B}^T\mathbf{y}$. Note that when $\lambda = 0$, we have the standard ordinary least squares estimate for $\hat{\boldsymbol{\alpha}}$.

Optimal selection of the penalisation parameter λ is a subtle and important matter. A value which is too small leads to 'overfitting', i.e. capturing the noise in the data. Conversely, a value which is too large leads to over smoothing of the data, i.e. 'underfitting'. Several criteria have been traditionally proposed (e.g. Hurvich et al. (1998), Wood (2006)) but the authors tackled this problem using a Bayesian modelling framework which is detailed in Evers et al. (2014).

References

- Ahuja, N., Schacter, B.J., 1983. *Pattern Models*. John Wiley & Sons, New York.
- Aziz, J.A., Newell, C.J., Ling, M., Rifai, H.S., Gonzales, J.R., 2003. Maros: a decision support system for optimizing monitoring plans. *Ground Water* 41 (3).
- Barcelona, M.J., Gibb, J.P., Helfrich, J.A., Garske, E.E., 1985. *Practical Guide for Ground-water Sampling*. SWS Contract Report 374. URL: <http://www.epa.gov/oust/cat/pracgw.pdf>.
- Bowman, A., Azzalini, A., 1997. *Applied Smoothing Techniques for Data Analysis: the Kernel Approach with S-plus Illustrations*. Oxford University Press, Oxford.
- Bowman, A., Crawford, E., Alexander, G., Bowman, R.W., 2007. rpanel: simple interactive controls for R functions using the tcltk package. *J. Stat. Softw.* 17 (9), 1–18. URL: <http://www.jstatsoft.org/v17/i09/>.
- Bowman, A.W., Azzalini, A., 2010. *R Package Sm: Nonparametric Smoothing Methods* (Version 2.2–4). University of Glasgow/Università di Padova, UK/Italia.
- URL: <http://www.stats.gla.ac.uk/~adrian/sm>. http://azzalini.stat.unipd.it/Book_sm.
- Cameron, K., 2004. Better optimization of ltm networks. *Bioremediat. J.* 8 (03–04).
- Cameron, K., Hunter, P., 2002. Using spatial models and kriging techniques to optimize long-term ground-water monitoring networks: a case study. *Environmetrics* 13, 629–656.
- Carslaw, D.C., Ropkins, K., 2012. openair — an R package for air quality data analysis. *Environ. Model. Softw.* 27–28 (0), 52–61.
- Cressie, N., Wikle, C.K., 2011. *Statistics for Spatio-temporal Data*. Wiley, New York.
- Eilers, P.H.C., Marx, B.D., 1992. Generalized linear models with P-splines. In: Fahrmeir, L., et al. (Eds.), *Advances in GLIM and Statistical Modelling*. Springer, New York.
- Eilers, P.H.C., Marx, B.D., 1996. Flexible smoothing with b-splines and penalties. *Stat. Sci.* 11, 89–121.
- Esri, 1998. *Esri Shapefile Technical Description*. URL: <http://www.esri.com/library/whitepapers/pdfs/shapefile.pdf>.
- Evers, L., Molinari, D.A., Bowman, A.W., Jones, W.R., Spence, M.J., 2014. Efficient and automatic methods for flexible regression on spatiotemporal data, with applications to groundwater monitoring. *Environmetrics* (in press).
- Gaus, I., Kinniburgh, D.G., Talbot, J.C., Webster, R., 2003. Geostatistical analysis of arsenic concentration in groundwater in Bangladesh using disjunctive kriging. *Environ. Geol.* 44 (8).
- Helsel, D.R., 2004. *Nondetects and Data Analysis*. John Wiley & Sons, New York.
- Helsel, D.R., Hirsch, R.M., 2002. *Statistical Methods in Water Resources*. Tech. Rep., United States Geological Survey. URL: <http://water.usgs.gov/pubs/twri/twri4a3/>.
- Hirsch, R.M., Slack, J.R., Smith, R.A., 1982. Techniques of trend analysis for monthly water-quality data. *Water Resour. Res.* 18 (1), 107–121.
- Hurvich, C., Simonoff, J., Tsai, C.-L., 1998. Smoothing parameter selection in nonparametric regression using an improved akaike information criterion. *J. R. Stat. Soc. Ser. B* 60, 271–293.
- Lang, D.T., 2012. RDCOMClient: R-DCOM Client. R Package Version 0.93-0. URL: <http://www.omegahat.org/RDCOMClient>.
- Lewin-Koh, N.J., Bivand, R., contributions by Edzer J., Pebesma, Archer, E., Baddeley, A., Bibiko, H.-J., Callahan, J., Carrillo, G., Dray, S., Forrest, D., Friendly, M., Giraudoux, P., Golicher, D., Rubio, V.G., Hausmann, P., Hufthammer, K.O., Jagger, T., Luque, S.P., MacQueen, D., Niccolai, A., Short, T., Snow, G., Stabler, B., Turner, R., 2012. *Maptools: Tools for Reading and Handling Spatial Objects*. R Package Version 0.8-16. URL: <http://CRAN.R-project.org/package=maptools>.
- McLeod, A., 2011. Kendall: Kendall Rank Correlation and Mann-Kendall Trend Test. R Package. URL: <http://www.stats.uwo.ca/faculty/aim>.
- Pebesma, E.J., Bivand, R.S., November 2005. Classes and methods for spatial data in R. *R. News* 5 (2), 9–13. URL: <http://CRAN.R-project.org/doc/Rnews/>.
- R Development Core Team, 2012. *R: a Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, ISBN 3-900051-07-0. URL: <http://www.R-project.org/>.
- Ricker, J.A., 2008. A practical method to evaluate ground water contaminant plume stability. *Ground Water Monit. Remediat.* 28 (4), 85–94.
- Rowlingson, B., Diggle, P., Adapted, Packaged for R by Roger Bivand, pcp Functions by Giovanni Petris, Goodness of Fit by Stephen Eglen, 2012. *Splancs: Spatial and Space-time Point Pattern Analysis*. R Package Version 2.01-31. URL: <http://CRAN.R-project.org/package=splancs>.
- Sarkar, D., 2008. *Lattice: Multivariate Data Visualization with R*. Springer, New York, ISBN 978-0-387-75968-5. URL: <http://lmdvr.r-forge.r-project.org>.
- Tierney, L., 2011. tkrplot: TK Rplot. R Package Version 0.0-23. URL: <http://CRAN.R-project.org/package=tkrplot>.
- Turner, R., 2012. deldir: Delaunay Triangulation and Dirichlet (Voronoi) Tessellation. R Package Version 0.0-19. URL: <http://CRAN.R-project.org/package=deldir>.
- Wood, S.N., 2006. *Generalized Additive Models — an Introduction with R*. Chapman & Hall/CRC.
- Xie, Y., 2012. *Animation: a Gallery of Animations in Statistics and Utilities to Create Animations*. R Package Version 2.1. URL: <http://CRAN.R-project.org/package=animation>.